

# **INFORME FINAL:**

## **Prediccion de Engagement en Facebook**

Analisis Comparativo, Diagnostico y Modelado Predictivo

*Dataset: Cosmetic Brand (Moro et al., 2016)*

Generado automaticamente tras analisis exhaustivo en Python

## 1. Introduccion y Contexto de Negocio

El presente informe detalla el proceso de analisis, limpieza y modelado predictivo realizado sobre un conjunto de datos perteneciente a una reconocida marca de cosmeticos en Facebook. El objetivo principal no es solo predecir metricas, sino entender la dinamica detras de la interaccion del usuario.

A diferencia de enfoques tradicionales, este proyecto distingue dos necesidades de negocio fundamentales:

- Necesidad A (Prediccion a Priori): Estimar el exito de un post ANTES de publicarlo, usando solo variables disponibles al momento de la creacion (Hora, Mes, Tipo).
- Necesidad B (Diagnostico Post-Hoc): Entender las causas del exito UNA VEZ publicado el post, analizando correlaciones con la viralidad (Shares, Reach).

## 2. Hallazgos del Analisis Exploratorio (EDA)

Se realizo un analisis estadistico riguroso que arrojo los siguientes descubrimientos clave:

- Distribucion del Target: La variable 'Likes' presenta una distribucion log-normal con fuerte sesgo positivo. El test de Shapiro-Wilk confirmo la no-normalidad ( $p < 0.05$ ), justificando la transformacion logaritmica (Log1p) para el modelado.
- Diferencia con la Literatura: Contrario al paper de Moro et al., en este dataset el 'Tipo de Contenido' (Foto vs Video) mostro menor relevancia que las variables temporales y el tamaño de la comunidad (Page Likes).
- Data Leakage: Se detecto una correlacion de Spearman  $> 0.85$  entre 'Shares' y 'Likes', lo que confirma que la viralidad es el predictor mas potente, aunque no este disponible 'a priori'.

## 3. Definicion de Escenarios de Modelado

Para garantizar la honestidad metodologica, se diseñaron tres escenarios experimentales:

### Escenario 1: Paper Original (Benchmark)

Replicacion exacta de las 7 variables usadas por Moro et al. para establecer una linea base.

### Escenario 2: Optimizado (Sin Leakage)

Ingenieria de caracteristicas temporal. Se crearon variables como 'Is\_Weekend' y 'Time\_Segment' (Morning/Evening) para mejorar la prediccion a priori sin usar datos del futuro.

### Escenario 3: Lifetime (Diagnostico)

Inclusion de metricas post-publicacion seleccionadas via Information Gain (Mutual Information). Este escenario modela la relacion Viralidad-Exito.

## 4. Resultados del Benchmark y Evaluacion

Se evaluaron 4 algoritmos (Ridge, RF, XGBoost, SVM) utilizando validacion cruzada de 5 pliegues (K-Fold). La metrica principal fue el MAPE (Error Porcentual Absoluto Medio).

Modelo	Esc 1 (Paper)	Esc 2 (Pred)	Esc 3 (Diag)
Ridge Regression	139.9%	136.2%	>1000% (Fail)
SVM (RBF)	137.8%	157.8%	45.7%
XGBoost	<b>105.5% (Best)</b>	105.8%	38.4%
Random Forest	115.0%	114.6%	<b>36.8% (Best)</b>

Analisis de Resultados:

- La predicción a priori (Esc 1 y 2) presenta un error alto (~105%), indicando que el éxito en redes sociales tiene un componente estocástico alto que no depende solo de la hora de publicación.
- El modelo de diagnóstico (Esc 3) reduce el error drásticamente al 36% con Random Forest, confirmando que el sistema es predecible si se conoce el alcance (Reach) y la viralidad (Shares).

## 5. Interpretacion (SHAP) y Conclusiones

El análisis de valores SHAP reveló la 'Caja Negra' de los modelos:

1. En el Escenario 3, la variable 'Shares' domina la predicción. Un post compartido actúa como multiplicador exponencial de Likes.
2. En el Escenario 2, las variables temporales como 'Month' y 'Hour' tienen impacto, pero insuficiente para explicar los outliers virales.

CONCLUSION FINAL:

El proyecto demuestra que es posible diagnosticar con precisión el rendimiento ( $R^2 > 0.6$ ) basándose en la viralidad. Sin embargo, la predicción 'a ciegas' requiere incorporar nuevas fuentes de datos (análisis de imagen/texto) para superar la barrera del 100% de error.